

Mark Zhao

myzhao@stanford.edu | 418 Gates Computer Science, Stanford, CA 94305 | (662)-801-1496
<https://web.stanford.edu/~myzhao/>

Research Interests

I build performant and scalable **systems for machine learning** (ML) to train, serve, and enable applications with large-scale machine learning models. My current research focuses on co-designing the **interacting components that compose modern ML systems**, including the ML training data pipeline and compound systems for ML inference. I am broadly interested in applying tools across cloud computing systems, machine learning, and computer architecture.

Education

Stanford University <i>Ph.D. in Electrical Engineering</i> Dissertation: <i>Performant and Scalable Systems Across the Machine Learning Pipeline</i> Advisor: Christos Kozyrakis	2025 (expected)
Cornell University <i>B.S. in Electrical and Computer Engineering, summa cum laude</i> Research Advisor: Edward Suh	2018

Industry Research Experience

Meta Platforms <i>Visiting Researcher, FAIR SysML & Capacity Engineering and Analysis</i> Mentors: Carole-Jean Wu and Niket Agarwal · Built, deployed, and optimized distributed systems to improve the performance and efficiency of Meta's production machine learning infrastructure. Projects included a disaggregated data preprocessing service (DPP), a flash storage tier for ML datasets (Tectonic-Shift), and deduplication optimizations for recommendation model training infrastructure (RecD).	2020 – 2022
Intel Corporation <i>Graduate Cloud Engineering Intern, Data Center Group</i> Mentor: Arindam Saha · Developed an inference serving framework that dynamically manages ML accelerator designs on Intel FPGAs to maximize serving performance across diverse inference requests.	2019

Peer-Reviewed Publications

cedar: Optimized and Unified Machine Learning Input Data Pipelines Mark Zhao , Emanuel Adamiak, and Christos Kozyrakis [VLDB 2025 (Accepted with Shepherding)] <i>Proceedings of the VLDB Endowment</i> , Volume 18	2025
---	------

- ReCycle: Resilient Training of Large DNNs using Pipeline Adaptation** 2024
Swapnil Gandhi, **Mark Zhao**, Athinagoras Skiadopoulos, and Christos Kozyrakis
[**SOSP 2024**] 30th Symposium on Operating Systems Principles
- High-throughput and Flexible Host Networking for Accelerated Computing** 2024
Athinagoras Skiadopoulos, Zhiqiang Xie, **Mark Zhao**, Qizhe Cai, Saksham Agarwal, Jacob Adelman, David Ahern, Carlo Contavalli, Michael Goldflam, Vitaly Mayatskikh, Raghu Raja, Daniel Walton, Rachit Agarwal, Shrijeet Mukherjee, and Christos Kozyrakis
[**OSDI 2024**] 2024 USENIX Symposium on Operating Systems Design and Implementation
- Tectonic-Shift: A Composite Storage Fabric for Large-Scale ML Training** 2023
Mark Zhao, Satadru Pan, Niket Agarwal, Zhaoduo Wen, David Xu, Anand Natarajan, Pavan Kumar, Shiva Shankar P, Ritesh Tijoriwala, Karan Asher, Hao Wu, Aarti Basant, Daniel Ford, Delia David, Nezhir Yigitbasi, Pratap Singh, Carole-Jean Wu, and Christos Kozyrakis
[**ATC 2023**] 2023 USENIX Annual Technical Conference
Invited fast-track submission to ACM Transactions on Storage
- RecD: Deduplication for End-to-End Deep Learning Recommendation Model Training Infrastructure** 2023
Mark Zhao, Dhruv Choudhary, Devashish Tyagi, Ajay Somani, Max Kaplan, Sung-Han Lin, Sarunya Pumma, Jongsoo Park, Aarti Basant, Niket Agarwal, Carole-Jean Wu, and Christos Kozyrakis
[**MLSys 2023**] 6th Conference on Machine Learning and Systems
- Understanding Data Storage and Ingestion for Large-Scale Deep Recommendation Model Training** 2022
Mark Zhao, Niket Agarwal, Aarti Basant, Buğra Gedik, Satadru Pan, Mustafa Ozdal, Rakesh Komuravelli, Jerry Pan, Tianshu Bao, Haowei Lu, Sundaram Narayanan, Jack Langman, Kevin Wilfong, Harsha Rastogi, Carole-Jean Wu, Christos Kozyrakis, and Parik Pol
[**ISCA 2022**] 49th Annual International Symposium on Computer Architecture
- ShEF: Shielded Enclaves for Cloud FPGAs** 2022
Mark Zhao, Mingyu Gao, and Christos Kozyrakis
[**ASPLOS 2022**] 27th ACM International Conference on Architectural Support for Programming Languages and Operating Systems
- Llama: A Heterogeneous & Serverless Framework for Auto-tuning Video Analytics Pipelines** 2021
Francisco Romero*, **Mark Zhao***, Neeraja J Yadwadkar, and Christos Kozyrakis
[**SoCC 2021**] 12th ACM Symposium on Cloud Computing
(* denotes equal contribution)
- HyperFlow: A High-Assurance Processor Architecture for Practical Timing-Safe Information Flow Security** 2018
Andrew Ferraiuolo, **Mark Zhao**, Andrew C. Myers, and G. Edward Suh
[**CCS 2018**] 25th ACM Conference on Computer and Communications Security
- FPGA-Based Remote Power Side-Channel Attacks** 2018
Mark Zhao and G. Edward Suh
[**S&P 2018**] 39th IEEE Symposium on Security and Privacy
Distinguished Practical Paper Award
2022 Top Pick in Hardware and Embedded Security

Technical Articles

- Remote Power Side-Channel Attacks on FPGAs** 2024
Mark Zhao and G. Edward Suh
IEEE Design & Test, 2024
- Counting Spree: Color Recognition and Segmentation in Real-time Video to Detect Manufacturing Defects** 2018
Mark Zhao and Claire Chen
Circuit Cellar Magazine, Issue #333, April 2018

Awards and Honors

- Meta Ph.D. Fellowship in AI System HW/SW Co-Design** 2023
· *Full funding and stipend for two academic years*
- MLCommons Machine Learning and Systems Rising Star** 2023
- Top Pick in Hardware and Embedded Security** 2022
· *For FPGA-based Remote Power Side-Channel Attacks*
- Stanford Graduate Fellowship** 2018
· *Full funding and stipend for three academic years*
- Distinguished Practical Paper Award**, IEEE Symposium on Security and Privacy 2018
· *For FPGA-based Remote Power Side-Channel Attacks*
- Sibley Prize**, Cornell ECE 2018
· *Awarded to the top graduating senior in Electrical and Computer Engineering*
- Meinig Family Cornell National Leadership Scholar**, Cornell University 2014
· *University-wide scholarship for demonstrating "an outstanding degree of leadership"*
- United States Presidential Scholar**, U.S. Department of Education 2014
· *Program established in 1964, by executive order of the President, to "recognize and honor some of our nation's most distinguished graduating high school seniors"*

Invited Talks

- End-to-End Optimization of Large-Scale ML Training Systems**
· *AMD Research and Advanced Development* 2024
· *UCF ECE Computer Architecture Seminar Series* 2024
· *SRC JUMP 2.0 ACE Center for Evolvable Computing Annual Review* 2023
- Understanding and Optimizing Data Storage and Ingestion Systems**
· *SRC JUMP 2.0 ACE Center for Evolvable Computing Liason Meeting* 2023
· *Cornell Systems Lunch* 2023
· *ByteDance Infrastructure Research Group* 2023
· *Stanford SystemX Fall Conference* 2022
- FPGA-Based Remote Power Side-Channel Attacks**
· *Top Picks in Hardware and Embedded Security Workshop* 2022

Llama: A Heterogeneous & Serverless Framework for Auto-Tuning Video Analytics Pipelines

- *Stanford Systems Seminar* 2021
- *Stanford Platform Lab Retreat* 2020

ShEF: Shielded Enclaves for Cloud FPGAs

- *Stanford SystemX Fall Conference* 2019
- *Stanford Platform Lab Review* 2019

Teaching Experience

CS 349D: Cloud Computing Technology , <i>Course Assistant</i> Stanford University	Spring 2024
CS 349D: Cloud Computing Technology , <i>Course Assistant</i> Stanford University	Spring 2023
EE 180: Digital Systems Architecture , <i>Course Assistant</i> Stanford University	Winter 2023
ECE 5760: Advanced Microcontroller Design , <i>Teaching Assistant</i> Cornell University	Spring 2018
ECE 4760: Designing with Microcontrollers , <i>Teaching Assistant</i> Cornell University	Fall 2017
ECE 3140: Embedded Systems , <i>Teaching Assistant</i> Cornell University	Spring 2017
PHYS 2213: Physics II (Electromagnetism) , <i>Undergraduate Teaching Assistant</i> Cornell University	Fall 2015
MATH 1920: Multivariable Calculus for Engineers , <i>Course Assistant</i> Cornell University	Fall 2015

Service

Stanford EE Faculty Search Committee , <i>Graduate Student Member</i>	2024
Workshop on ML for Computer Architecture and Systems (MLArchSys at ISCA'24) , <i>Technical Program Committee</i>	2024
Workshop on Machine Learning and Systems (EuroMLSys at EuroSys'24) , <i>Technical Program Committee</i>	2024
Workshop on ML for Computer Architecture and Systems / Architecture and System Support for Transformer Models Workshop (MLArchSys/ASSYST at ISCA 2023) , <i>Technical Program Committee</i>	2023
IEEE Transactions on Circuits and Systems II: Express Briefs , <i>External Reviewer</i>	2022
Design Automation Conference (DAC) , <i>External Reviewer</i>	2019

Mentorship

Zhanqiu (Summer) Hu (Ph.D. @ Cornell Tech) · <i>End-to-End Optimization of Recommendation Systems</i>	2024– Present
---	------------------

Suze van Adrichem (B.S. @ Stanford) · <i>PandoRT: A Distributed Serving System for Compound LLM Applications</i>	2024– Present
Jenny Wei (B.S. @ Stanford) · <i>Building and Optimizing Systems for LLM Pipelines</i>	2024– Present
Laasya Konidala (B.S. @ Stanford) · <i>Building and Optimizing RAG for LLM Pipeline Serving Systems</i>	2024
Ethan Zhang (B.S. @ Stanford) · <i>A New Frontier for Model Routing</i>	2024
Emanuel Adamiak (B.S. @ Stanford) · <i>cedar: Optimized and Unified Machine Learning Input Data Pipelines</i>	2023 – 2024
Andrew Woen (B.S. @ Stanford) · <i>Optimizing Data Storage and Ingestion Pipelines for ML Training</i>	2023